

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平7-92988

(43) 公開日 平成7年(1995)4月7日

(51) IntCl. ⁸	識別記号	庁内整理番号	F I	技術表示箇所
G 1 0 L 3/00	5 1 1	9379-5H		
H 0 4 R 3/00	3 2 0			

審査請求 未請求 請求項の数27 O L (全 14 頁)

(21) 出願番号 特願平5-238579

(22) 出願日 平成5年(1993)9月27日

(71) 出願人 000005821

松下電器産業株式会社

大阪府門真市大字門真1006番地

(72) 発明者 則松 武志

大阪府門真市大字門真1006番地 松下電器
産業株式会社内

(72) 発明者 中藤 良久

大阪府門真市大字門真1006番地 松下電器
産業株式会社内

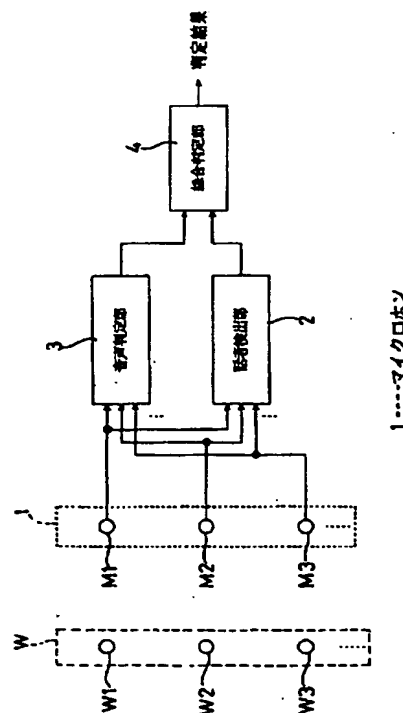
(74) 代理人 弁理士 森本 義弘

(54) 【発明の名称】 音声検出装置と映像切り替え装置

(57) 【要約】

【目的】 話者の音声が発見でき話者に対応したマイクロホン1を正確に特定できる音声検出装置と、この特定により映像を自動的に話者に切り換えることができる映像切り替え装置を提供することを目的とする。

【構成】 音声判定部3が、マイクロホン1に入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定する。話者検出部2が、隣接したマイクロホン1の入力信号の間の差異を検出することにより話者の位置を推定し、この話者に対応したマイクロホン1を特定する。以上の音声判定部3と話者検出部2の出力結果に基づいて、総合判定部4がそれぞれのマイクロホン1に対応した話者の音声のみを判定する。



【特許請求の範囲】

【請求項 1】 音響を検出する複数のマイクロホンと、これらのマイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定する音声判定部と、任意のマイクロホンの入力信号とこのマイクロホンに隣接した位置にあるマイクロホンの入力信号との間の差異を検出することにより音響の発生源である話者の位置を推定し、この話者に対応したマイクロホンを特定する話者検出部と、前記音声判定部と話者検出部の出力結果を用いて予め定めた判定条件をもとにそれぞれのマイクロホンに対応した話者の音声のみを判定する総合判定部とを備えた音声検出装置。

【請求項 2】 話者検出部を、隣接する 2 つのマイクロホンの入力信号間の相互相関係数を用いて隣接する前記マイクロホンへの入力信号の到達時間の差を検出することにより、話者の位置を推定し、この話者に対応したマイクロホンを特定するよう構成した請求項 1 に記載の音声検出装置。

【請求項 3】 話者方向に向いた第 1 のマイクロホンと、話者と反対方向に向いた第 2 のマイクロホンと、前記第 1 のマイクロホンと第 2 のマイクロホンのそれぞれの入力信号の差異を検出することにより第 1 のマイクロホンの前方より発せられた信号のみを検出する前方音検出部と、第 1 のマイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定する音声判定部と、前記前方音検出部と音声判定部の出力結果を用いて予め定めた判定条件をもとにそれぞれの第 1 のマイクロホンに対応した話者の音声のみを判定する総合判定部とを備えた音声検出装置。

【請求項 4】 話者方向に向いた第 1 のマイクロホンと話者と反対方向に向いた第 2 のマイクロホンとを一組とする複数組のマイクロホンと、それぞれの組の前記第 1 のマイクロホンと第 2 のマイクロホンのそれぞれの入力信号の差異を検出することにより第 1 のマイクロホンの前方より発せられた信号のみを検出する前方音検出部と、それぞれの組の第 1 のマイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定する音声判定部と、任意の第 1 のマイクロホンの入力信号とこのマイクロホンに隣接した位置にある第 1 のマイクロホンの入力信号との間の差異を検出することにより話者の位置を推定し、この話者に対応したマイクロホンを特定する話者検出部と、前記前方音検出部と音声判定部及び話者検出部の出力結果を用いて予め定めた判定条件をもとにそれぞれの組の第 1 のマイクロホンに対応した話者の音声のみを判定する総合判定部とを備えた音声検出装置。

【請求項 5】 前方音検出部を、第 1 のマイクロホンと

第 2 のマイクロホンのそれぞれの入力信号のパワーの差を算出し、この値により第 1 のマイクロホンの前方より発せられた信号であるか否かを判定するよう構成した請求項 3 または請求項 4 のいずれかに記載の音声検出装置。

【請求項 6】 前方音検出部を、第 1 のマイクロホンと第 2 のマイクロホンのそれぞれの入力信号のパワーの比を算出し、この値により第 1 のマイクロホンの前方より発せられた信号であるか否かを判定するよう構成した請求項 3 または請求項 4 のいずれかに記載の音声検出装置。

【請求項 7】 話者検出部を、隣接する 2 つの第 1 のマイクロホンの入力信号間の相互相関係数を用いて隣接する前記第 1 のマイクロホンへの入力信号の到達時間の差を検出することにより、話者の位置を推定し、この話者に対応した第 1 のマイクロホンを特定するよう構成した請求項 4 に記載の音声検出装置。

【請求項 8】 音声判定部を、予め多数の音声データから音声信号の持つ周波数的特徴あるいは時間的特徴を求めておき、入力信号がどの程度前記周波数的特徴あるいは時間的特徴が類似しているかを表す指標により音声と雑音を判別し、前記周波数的特徴あるいは時間的特徴を持つ音声信号のみを検出するよう構成した請求項 1 または請求項 3 または請求項 4 のいずれかに記載の音声検出装置。

【請求項 9】 音声判定部を、入力信号を線形予測分析した際に得られた線形予測係数あるいはケプストラム係数あるいは自己相関係数を、予め作成しておいた音声に関する前記線形予測係数あるいはケプストラム係数あるいは自己相関係数と比較することにより周波数的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 10】 音声判定部を、予め作成しておいた音韻毎のスペクトルと入力信号のスペクトルがどの程度似通っているかに基づいて音声の音韻性を認識することにより周波数的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 11】 音声判定部を、周波数軸をデジタルあるいはアナログフィルタにより数帯域に分割し、前記デジタルあるいはアナログフィルタにより得られた各帯域毎のエネルギーのパターンを認識することにより周波数的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 12】 音声判定部を、周波数軸をデジタルあるいはアナログフィルタにより数帯域に分割し、前記デジタルあるいはアナログフィルタにより得られた各帯域毎の信号の零交差を求め、各帯域毎の前記零交差の回数

により周波数的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 13】 音声判定部を、周波数軸をデジタルあるいはアナログフィルタにより数帯域に分割し、前記デジタルあるいはアナログフィルタにより得られた各帯域毎の信号の 1 次以上の自己相関係数により周波数的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 14】 音声判定部を、周波数軸をデジタルあるいはアナログフィルタにより数帯域に分割し、前記デジタルあるいはアナログフィルタにより得られた各帯域毎の信号を F F T 分析した際に得られた 1 次以上のケプストラム係数により周波数的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 15】 音声判定部を、周波数軸をデジタルあるいはアナログフィルタにより数帯域に分割し、前記デジタルあるいはアナログフィルタにより得られた各帯域毎の信号を F F T 分析した際に得られた 1 次以上の自己相関係数及び 1 次以上のケプストラム係数のうち少なくとも 1 つ以上の特徴量により周波数的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 16】 音声判定部を、周波数軸をデジタルあるいはアナログフィルタにより数帯域に分割し、前記デジタルあるいはアナログフィルタにより得られた各帯域毎の信号を F F T 分析し得られた特徴量をベクトル量子化して求めたコードブックにより周波数的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 17】 音声判定部を、話者の発声した音声中の音韻性の特徴付ける特徴量をベクトル量子化して求めたコードブックを予め求めておき、入力信号を前記コードブックにてベクトル量子化した際の量子化歪みにより周波数的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 18】 音声判定部を、入力信号のスペクトルが時事刻々いかなる変化をしているかに基づいて音声中の音韻性を認識することにより時間的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 19】 音声判定部を、予め多数の音声から求めておいた音韻毎の継続時間の最大値および最小値により入力信号から分析フレーム毎に音韻を検出し、各音韻

がどの程度継続しているかを示す継続時間を求め、前記音韻毎の継続時間の最大値より小さくしかも最小値より大きいときのみ音声が入力されたとすることにより時間的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 20】 音声判定部を、予め多数の音声から求めておいた音韻毎のスペクトル系列を標準モデルとして予め求めておき、前記標準モデルを用いて入力信号中のスペクトルがどの程度継続しているかを表す継続時間を計測することにより時間的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 21】 音声判定部を、話者の発声した音声中の音韻性の特徴付ける特徴量をベクトル量子化して求めたコードブックを用いて、入力信号をベクトル量子化した際のコード列の変化のパターンを認識することにより時間的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 22】 音声判定部を、話者の発声した音声中の音韻性の特徴付ける特徴量をベクトル量子化して求めたコードブックを用いて、入力信号をベクトル量子化した各コードがどの程度継続して現れるかにより時間的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 23】 音声判定部を、予め多数の音声データから各音韻毎の HMM モデルを作成しておき、前記 HMM モデルを用いて入力信号中に存在する音韻性を認識することにより周波数的特徴あるいは時間的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 24】 音声判定部を、入力信号から分析フレーム毎に音声を特徴付ける特徴量を抽出し、入力信号中の音声成分がどの程度継続しているか予め多数の音声データより求めておいた継続時間に関するファジィメンバーシップ関数を用いてファジィ推論することにより時間的特徴を検出して、入力信号の音声と雑音を判別し、入力信号中の音声のみを検出するよう構成した請求項 8 に記載の音声検出装置。

【請求項 25】 請求項 1 に記載の音声検出装置と、各話者の映像を出力するために、それぞれの話者の位置を予め記憶し出力映像を制御するカメラ制御部と、前記音声検出部の出力に基づいて音声が入力されているマイクロホンを選定し、対応する話者の映像に切り換えるための制御信号を前記カメラ制御部に出力する映像切り替え制御部とを備えた映像切り替え装置。

【請求項 26】 請求項 3 に記載の音声検出装置と、各

話者の映像を出力するために、それぞれの話者の位置を予め記憶し出力映像を制御するカメラ制御部と、前記音声検出部の出力に基づいて音声が入力されている第1のマイクロホンとを特定し、対応する話者の映像に切り換えるための制御信号を前記カメラ制御部に出力する映像切り替え制御部とを備えた映像切り替え装置。

【請求項27】 請求項4に記載の音声検出装置と、各話者の映像を出力するために、それぞれの話者の位置を予め記憶し出力映像を制御するカメラ制御部と、前記音声検出部の出力に基づいて音声が入力されている第1のマイクロホンとを特定し、対応する話者の映像に切り換えるための制御信号を前記カメラ制御部に出力する映像切り替え制御部とを備えた映像切り替え装置。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、テレビ会議システム等における話者の位置を特定する音声検出装置とこの出力により映像を切り替える映像切り替え装置に関するものである。

【0002】

【従来の技術】 近年、ISDN等デジタル通信網の発達により、企業の間では遠隔地間で積極的にテレビ会議システムを利用し始めている。

【0003】 現在のテレビ会議システムにおいて、限られた大きさのモニター画面を用いてより自然な会議進行を実現するためには、発言者が誰であるのかを知らせるためにリアルタイムにモニター画面を発言者に切り換える必要がある。現在の多くの会議システムでは、発言者が切り替わる度に操作卓を使ってマニュアルで映像を切り換えなければならず、自然な会議の進行の妨げになっていた。そこで会議中の発言者の音声を自動的に検出し発言者の映像に自動的に切り換えるための音声検出装置の実現が望まれている。

【0004】 実際に複数の参加者が存在するテレビ会議の場面を想定すると、会議中には参加者の発言した音声以外に様々な雑音が発生する。また全参加者の音声を收音するために会議室には複数のマイクロホンが設置されることになるが、ある話者の音声は自分自信のマイクロホンだけでなく隣接した位置にあるマイクロホンにも入力される。さらに会議の相手方の音声も漏れ聞こえ各マイクロホンに混入する。このような状況下で上記の音声検出装置を実現するためには、入力信号から音声信号の部分を正確に判別すると共に、どのマイクロホンに対応した位置にいる話者の発声した音声であることを的確に判定できなければならない。

【0005】 このような音声検出装置を実現するために、各マイクロホンに入力される信号のパワーを算出し、パワーが検出されたときにそのマイクロホンに音声が入力されていると判断することによって、予め記憶されたそのマイクロホンに対応する話者の位置へ自動的に

カメラを向け映像を切り換える試みが行われている。ここでパワーが検出された区間が一定時間以下の場合は音声と判定しないことで突発的な雑音による誤判定を防止している。またある話者の音声と同時に隣接した複数のマイクロホンに混入し、複数のマイクロホン入力が音声であると判定される場合に対応するため、パワー強度の大きい方を選択する方法もある。

【0006】

【発明が解決しようとする課題】 しかしながら上記の構成では、突発的な雑音は取り除けるが、パワーの大きな連続的な信号であれば音声あるいは雑音にかかわらず反応してしまい、発言していない話者に誤って映像が切り替わる場合が発生するという問題点がある。

【0007】 また、発言者は必ずしもマイクロホンの正面から発声するとは限らず、口元とマイクロホンとの位置関係は変化するため、パワー強度の違いだけでは、どの話者の発声した音声であるかは正確には判定することができないという問題点もある。

【0008】 本発明は、上記従来の課題を解決するものであり、入力された信号が突発的、連続的なものにかかわらず正確に音声信号であるか否かが判別できる共に、その音声信号がそれぞれのマイクロホンに対応した話者から発声されたものであるかが正確に判定することができる音声検出装置と、この音声検出装置の判定結果に基づいて自動的に話者の映像を切り換えることができる映像切り替え装置を提供することを目的とする。

【0009】

【課題を解決するための手段】 請求項1に記載の音声検出装置は、音響を検出する複数のマイクロホンと、これらのマイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定する音声判定部と、任意のマイクロホンの入力信号とこのマイクロホンに隣接した位置にあるマイクロホンの入力信号との間の差異を検出することにより音響の発生源である話者の位置を推定し、この話者に対応したマイクロホンを特定する話者検出部と、前記音声判定部と話者検出部の出力結果を用いて予め定めた判定条件をもとにそれぞれのマイクロホンに対応した話者の音声のみを判定する総合判定部とを備えたことを特徴とする。

【0010】 請求項3に記載の音声検出装置は、話者方向に向いた第1のマイクロホンと、話者と反対方向に向いた第2のマイクロホンと、前記第1のマイクロホンと第2のマイクロホンのそれぞれの入力信号の差異を検出することにより第1のマイクロホンの前方より発せられた信号のみを検出する前方音検出部と、第1のマイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定する音声判定部と、前記前方音検出部と音声判定部の出力結果を用いてそれぞ

れの第1のマイクロホンに対応した話者の音声のみを判定する総合判定部とを備えたことを特徴とする。

【0011】請求項4に記載の音声検出装置は、話者方向に向いた第1のマイクロホンと話者と反対方向に向いた第2のマイクロホンとを一組とする複数組のマイクロホンと、それぞれの組の前記第1のマイクロホンと第2のマイクロホンのそれぞれの入力信号の差異を検出することにより第1のマイクロホンの前方より発せられた信号のみを検出する前方音検出部と、それぞれの組の第1のマイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定する音声判定部と、任意の第1のマイクロホンの入力信号とこのマイクロホンに隣接した位置にある第1のマイクロホンの入力信号との間の差異を検出することにより話者の位置を推定し、この話者に対応したマイクロホンを特定する話者検出部と、前記前方音検出部と音声判定部及び話者検出部の出力結果を用いて予め定めた判定条件をもとにそれぞれの組の第1のマイクロホンに対応した話者の音声のみを判定する総合判定部とを備えたことを特徴とする。

【0012】請求項25に記載の映像切り替え装置は、請求項1に記載の音声検出装置と、各話者の映像を出力するために、それぞれの話者の位置を予め記憶し出力映像を制御するカメラ制御部と、前記音声検出部の出力に基づいて音声が入力されているマイクロホンを特定し、対応する話者の映像に切り換えるための制御信号を前記カメラ制御部に出力する映像切り替え制御部とを備えたことを特徴とする。

【0013】

【作用】請求項1の構成によると、音声判定部が、マイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定する。話者検出部が、隣接したマイクロホンの入力信号の間の差異を検出することにより話者の位置を推定し、この話者に対応したマイクロホンを特定する。以上の音声判定部と話者検出部の出力結果に基づいて、総合判定部がそれぞれのマイクロホンに対応した話者の音声のみを判定する。

【0014】請求項3の構成によると、前方音検出部が、話者方向に向いた第1のマイクロホンと話者と反対方向に向いた第2のマイクロホンに入力された信号の差異を検出して、第1のマイクロホンの前方より発せられた信号のみを検出する。音声判定部が、第1のマイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定する。以上の前方音検出部と音声判定部の出力結果に基づいて、総合判定部がそれぞれの第1のマイクロホンに対応した話者の音声のみを判定する。

【0015】請求項4の構成によると、前方音検出部が、一組にされた話者方向に向いた第1のマイクロホンと話者と反対方向に向いた第2のマイクロホンに入力された信号の差異を検出して、第1のマイクロホンの前方より発せられた信号のみを検出する。音声判定部が、各組の第1のマイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定する。話者検出部が、隣接した第1のマイクロホンの入力信号の間の差異を検出することにより話者の位置を推定し、この話者に対応したマイクロホンを特定する。以上の前方音検出部と音声判定部と話者検出部の出力結果に基づいて、総合判定部が各組の第1のマイクロホンに対応した話者の音声のみを判定する。

【0016】請求項25の構成によると、請求項1に記載の音声検出装置の出力に基づいて、映像切り替え制御部が、特定したマイクロホンに対応した話者に映像を切り換える制御信号をカメラ制御部に出力する。この制御信号により、カメラ制御部は予め記憶した話者の位置情報に基づいて出力映像の切り替えを制御する。

【0017】

【実施例】以下、本発明の音声検出装置の第1の実施例について図面を参照しながら説明する。

【0018】図1は本実施例の構成を示すブロック図である。図1において、Wは音声を発する話者、1はマイクロホン、2は隣接したマイクロホンの入力信号間の波形上の類似性を調べることにより話者の位置を推定する話者検出部、3は各マイクロホンの入力信号から音韻の特徴を抽出し、音声信号であるか否かを判定する音声判定部、4は音声判定部および話者検出部の結果をもとに、それぞれのマイクロホンに対してそれぞれの前方に位置する話者の音声信号が入力されているかを否かを判定し、この判定結果を出力する総合判定部である。

【0019】以下、上記音声検出装置の動作を説明する。ここでは一般的なテレビ会議の場面を想定し、話者が横一線に並んでいるとし、また各話者にそれぞれマイクロホンが設置されているものとする。

【0020】まず、マイクロホン1に入力された音響信号はアナログ／デジタル変換され、話者検出部2、音声判定部3にそれぞれ入力される。話者検出部2では隣合うマイクロホン同志での入力信号間の相関関係を調べることで話者の位置を推定する。ここで例えば話者W2が発言している場合を考える。話者W2の発声した音声はマイクロホンM2はもちろんその隣のマイクロホンM1、M3にも入力される（その他のマイクロホンにも入力されるがそのパワーは小さくなる）。また話者W2は常にマイクロホンM2の正面方向に在るわけではなく、話者W1、あるいは話者W3の方向に寄って発声しているかもしれない。これらの位置関係を示したのが図2である。もし話者がマイクロホンM2、M3から等距

離の地点 x にいるときは、音声信号の各マイクロホンへの到達時間は等しいが、話者が左右にずれることによって到達時間に差が生じる。そこでこの到達時間の差を検出することにより、話者のおおよその位置を推定することが可能となる。

【0021】図3は話者検出部2の動作を示す要部フローチャートである。以下図3のフローチャートに沿って*

$$h_i = \sum_{t=1}^n (b_t \cdot C_{t+i})$$

ただし $-m \leq i \leq m, m > 0$

【0023】ここで b_t 、 C_t は任意の時刻 t におけるサンプル値、 n は1フレームのサンプル数、 m は話者の左右のずれを検出するために予め設定された値であり、分析条件、マイクロホンと話者の位置関係により多少変わってくる。次にステップ32で、各マイクロホンの組毎に得られたそれぞれの $-m$ 次から m 次までの相互相関係数のうち最大値を与える相関係数の値及びその次数を記憶する。ステップ33では、各マイクロホンの組毎の相互相関係数の最大値の中から最大値を与えるマイクロホンの組を選択する。次にステップ34で、選択されたマイクロホンの組の最大相関係値を与える次数から話者の左右へのずれ幅を推定し、話者が対応するマイクロホンの正面方向に存在するか否かを判定する。例えば図2において話者W2の位置から発声された音声信号のマイクロホンM2、マイクロホンM3への到達時間の差 T は音の速度を c 、話者W2からマイクロホンM2までの距離 l 、マイクロホンM3までの距離 k として式2で表される。

【0024】

【数2】

$$T = \frac{|l - k|}{c} \quad \text{-----式2}$$

【0025】ここで最大相関係値を与える次数が $m1$ であったとすると、 T は $T_s \times m1$ (秒) に相当し、話者W2は地点 x からほぼこの時間に相当する距離分だけ左にすることがわかる。 T_s はサンプリング周期である。そこで予めマイクロホン正面方向の話者の音声を捉えるべき範囲を設定しておき、検出の結果その範囲内であれば話者が存在すると判定する。またマイクロホンM2及びM3からはほぼ等距離の地点 x を含む線上の近傍に音源が存在する場合は、特に入力されているマイクロホンは特定しないようにする。

【0026】最後にステップ35で、判定結果として、話者が発声していると特定されたマイクロホンについてはオン信号を、特定されなかったマイクロホンについてはオフの信号を送出する。ここで誤判定、及び短い発

*説明する。図3のステップ31で、まず隣合う2つのマイクロホンそれぞれの組について入力信号の相互相関係数を一定時間間隔毎 (以下フレームと呼ぶ) に式1により算出する。

【0022】

【数1】

-----式1

言、突発的な雑音による判定結果の短時間での切り替わりを防止するため、同一の判定結果が一定フレーム続いた場合に判定結果をオンにし、またマイクロホンの特定が一つもできない状態が一定フレーム以上続いたときにオフにするよう制御する。以上が話者検出部2の動作説明である。

20 【0027】次に音声判定部3の動作について説明する。図4は音声判定部3に関するブロック構成図である。図4において41は音声検出のための複数の特徴量を抽出する特徴抽出部で、1フレーム毎の特徴量を算出する。これらの特徴量は音声を検出するために用いられるものであり、音声に特有の性質を有している。本実施例では1次以上のケプストラム係数を用いる。他の特徴量としてたとえば線形予測分析の際に得られる自己相関係数や線形予測係数、PARCOR係数、メルケプストラム係数等を用いても差し支えない。あるいは他の音声分析、たとえばFFT分析により得られるスペクトル情報を用いても、音声の特徴を捉えていることでは同じであるので使用可能である。また、入力信号をアナログフィルタあるいはディジタルフィルタにより周波数軸上で数個の帯域に分割し、各帯域のエネルギーを算出してそれをひとつの特徴量として扱うこともできる。また各帯域毎に求めた零交差回数の特徴量として使用することや、各帯域毎にFFT分析して得られるメルケプストラム係数をひとつの特徴量として扱う、また各帯域毎にLPC分析により得られるスペクトルをひとつの特徴量として扱うことも可能である。

40 【0028】次に、42は予め信頼性の高い多数の学習用音声データについて特徴抽出部41で抽出した特徴量を用いて、音声の周波数的な標準パターンを作成する周波数パターン作成部である。標準パターンとしては、予め多数の音声データからスペクトルに関する特徴量を抽出しておき、各音韻毎にその特徴量を用いて標準パターンを作成する。本実施例では標準パターンとしては、特徴量の分布を多次元正規分布としたときの平均、共分散を用い、これを音韻毎に作成しておく。また他の分布として、たとえばガンマ分布やポアソン分布等を用いて

も差し支えない。さらにこの標準パターンとしては、学習用音声データを音韻毎に分類した後各音韻毎に作成した最適な標準パターンを用いたり、学習用音声データをベクトル量子化によりクラスタリングすることにより得られたコードを用いても、より精度の高い判定が可能となる。

【0029】43は特徴抽出部41から出力される入力信号のフレーム毎のケプストラム係数について周波数パターン作成部42にて作成した音韻毎の特徴量分布との距離すなわち尤度を計算し、ある閾値と比較することで音声であるかそれ以外かを判定する尤度判定部である。

【0030】44は予め信頼性の高い多数の学習用音声データから作成した音声の時間的な特徴を表現する時間パターンを作成する時間パターン作成部である。本実施例においては、多数の学習用音声データから作成した、音韻毎の継続時間に関する最大値、最小値を用いる。また、他の例として、継続時間分布たとえば正規分布やガンマ分布、ポアソン分布等を用いても差し支えない。

【0031】45は、尤度判定部43にて入力信号のうち音声と判定された部分について、時間パターン作成部44にて作成した時間パターンとを比較することで、入力信号が音声であったかそれ以外であったかを判定する最終判定部である。本実施例では、入力信号から各音韻がどの程度継続しているかを示す継続時間を求め、予め多数の音声から求めておいた音声の継続時間の最大値および最小値を用いて、最大値より小さくしかも前記最小値より大きいときのみ音声を検出されたとする。ここで、音声の継続時間の最大値および最小値にかえて、継続時間が統計的な分布特性を持つと仮定し、入力信号から得られた音声の継続時間をもとに確率を求め、その確率がある閾値より大きければ音声であると断定することも可能である。また、時間パターンとして多数の音声データから標準的な音声のスペクトル系列を標準パターンとして登録しておき、入力信号とこの標準パターンとの非線形伸縮（DPマッチング）により、入力信号のどの部分に各標準パターンが存在するかを検出（スポッティング）することで、音声であるかそれ以外かを判定することが可能である。また、時間パターンとして多数の音声スペクトル系列から隠れマルコフモデル（HMM）を予め標準パターンとして作成しておき、入力信号とこのHMMモデルとの確率計算により、入力信号のどの部分に各標準パターンが存在するかを検出（スポッティング）し、音声であるかどうかを判定することも可能である。また、時間パターンを用いて音声を検出するのではなく、入力信号を音声分析して得られた特徴量の変化量を時々刻々求め、その変化量を閾値判定することで音声中の音韻を検出し、音声と雑音を判別することも可能で

ある。さらに話者の発声した音声の中の音韻性を特徴付ける特徴量や、フィルタリング処理により各帯域毎に音声分析して得られた特徴量をベクトル量子化して求めたコードブックを用いて、入力信号をベクトル量子化した際の量子化歪みを閾値判定することで音声であるか雑音であるかを判定したり、さらに入力信号をベクトル量子化した際のコード列の変化のパターンに変換し、その各コードの出現頻度や、各コードの継続時間により、音声であるかどうかを判定することも可能である。

10 【0032】以下、音声判定部3の動作について図4のブロック構成図を参照しながら詳細に説明する。音響信号がマイクロホンを通して入力されると、特徴抽出部41でまず複数の特徴量が抽出される。本実施例ではケプストラム係数を用いて判定する。一定時間毎にK次の自己相関係数 $A_i(k)$ が算出され、さらに $A_i(k)$ は0次の自己相関係数 $A_i(0)$ で正規化される。ここで一定の時間間隔は、例えばサンプリング周波数を10KHzとして、200点（20ms）とし、この時間単位をフレームと呼ぶ。フレームiでのL次のケプストラム係数 $C_i(l)$ を線形予測分析により求める。ここでは、これらの特徴量が互いに独立であるとして、一括して1つのベクトル（m次元） x として扱うことにする。

【0033】周波数パターン作成部42では、予め多数の学習用音声データを用いて、各音韻毎に特徴抽出部41で得られる特徴量を抽出し、各音韻毎の周波数パターンを作成する。音韻としては母音や無声摩擦音、鼻音、有声破裂音、破擦音、流音、半母音等が考えられる。ここでは次の方法により音韻毎の平均値 μ_{ic} と共分散行列 Σ_{ic} を周波数パターンとして使用する。ただし、kは音韻番号、cは特徴量分布作成部にて得られた値であることを示し、 μ_{ic} はm次元のベクトル、 Σ_{ic} は $m \times m$ 次元のマトリックスである。学習用音韻データとしては、例えばある標準話者の音韻kの部分の学習用データから切り出して用いればよい。また、複数の話者の音声データを用いることで、話者の発声の変動に強い標準モデルを作成することができる。

【0034】尤度判定部43は、特徴抽出部41から出力されるフレーム毎の入力信号のいくつかの特徴量について、周波数パターン作成部42にて作成した各音韻毎の標準パターンと対数尤度を計算する部分である。ここで対数尤度とは、各特徴量の分布を多次元正規分布と仮定した場合の統計的距離尺度であり、ある音韻の標準パターンkに対するiフレーム目の入力ベクトル x_i の特徴量尤度 L_{ik} は、式3により計算される。

【0035】

【数3】

$$L_{ik} = -\frac{1}{2} \cdot (x_i - \mu_k)^t \cdot \Sigma_k^{-1} (x_i - \mu_k) - \frac{1}{2} \cdot \ln |\Sigma_k| + C$$

-----式3

【0036】ただし、 x_i は m 次元のベクトル (m 次元の特徴量) であり、 t は転置、 -1 は逆行列を示す。そして式4により、各音韻毎の対数尤度と予め決めておいた各音韻毎との閾値とを比較することで音韻の検出を行*10

*う。

【0037】

【数4】

$$L_{ik} = \frac{1}{2 \cdot N + 1} \cdot \sum_{i=-N}^N L_{ik} \geq L_{kTH} \quad \text{-----式4}$$

【0038】ただし、 L_{kTH} は各音韻 k に関する判定閾値 (対数尤度の閾値) である。時間パターン作成部44では、予め多数の学習用音声データを用いて、各音韻毎の継続時間の最大値 D_{max} 、最小値 D_{min} を求め、最終判定部45において、最終的な音声かそれ以外の雑音であるかの判定を行う。まず尤度判定部43にて検出された音韻の情報を最終判定部45に送り、各音韻が何フレーム継続したかすなわち各音韻毎の継続時間 D_k を求める。そして、この継続時間 D_k と時間パターン作成部43にて求めておいた各音韻毎の継続時間の最大値より大きくかつ最小値より小さいとき音韻が検出されたと判定し、最終的に入力信号が音声であるかそれ以外であるかを判定する。

【0039】さらに、このような音韻がある区間内でのくらいの頻度で出現するかを、ファジィ推論により判定することもできる。たとえば予め多数の音声データから各音韻毎の出現数に関するメンバシップ関数を決定しておき、実際に入力信号の各音韻毎の出現数を上記音韻判定部43にて求め、メンバシップ関数から算出されるファジィ出力を最終的に判定することで音声検出されたのか雑音検出されたのかを決定することができる。以上が音声判定部3の動作説明である。

【0040】最後に総合判定部4では、話者検出部2において対応する話者が発言しているとして特定されたマイクロホンの入力について、音声判定部3で音声信号が入力されていると判定されている場合に、そのマイクロホンはオンであるという信号を外部に送出する。

【0041】以上のように本実施例によれば隣接マイクロホン間の相関関係から話者方向から信号が入力されているマイクロホンを特定し、また音韻性を用いて入力信号が音声か否かを正確に判別することにより、突発雑音、連続的な雑音が入力されたときに誤って音声と誤判定するのを防ぐことができ、また音声信号が隣接するマイクロホンへ入力された場合でも話者に対応するマイクロホンを特定することができ、さらに周囲騒音等による誤反応をも防止することができる。

【0042】次に本発明の音声検出装置の第2の実施例について図面を参照しながら説明する。図5は第2の実施例の音声検出装置の構成を示すブロック図である。図5において、 W は音声を発する話者 (例えば、話者 $W1$ 、 $W2$ などで構成されている)、51は話者方向に向いた第1のマイクロホン (例えば、マイクロホン $M11$ 、 $M21$ などで構成されている)、52は話者と反対方向の向いた第2のマイクロホン (例えば、マイクロホン $M12$ 、 $M22$ などで構成されている)、53はマイクロホン51とマイクロホン52の入力信号から話者方向からの信号のみを検出する前方音検出部、54は第1のマイクロホンの入力信号からスペクトルの特徴量を検出し、音声であるか否かを判定する音声判定部、55は上記結果から話者方向からの音声信号のみを判定し、この判定結果を出力する最終判定部である。

【0043】以下、上記音声検出装置の動作を説明する。音響信号が各第1のマイクロホン51、第2のマイクロホン52に入力され、両方の信号が前方音検出部53に、第1のマイクロホンへの入力信号のみが音声判定部54に送出される。ここでは話者毎に第1のマイクロホンと第2のマイクロホンが一组として設置されているものとする。

【0044】前方音検出部53ではマイクロホン51、52のそれぞれの入力信号の差によりマイクロホン51の前方からの信号であるか否かを判定する。また、どの話者からの音声であるかの推定は、前方音検出部53によりマイクロホン51とマイクロホン52のそれぞれの入力信号のパワーの差を求め、この差が最も大きな値となるマイクロホン51の前方の話者からの音声であると判定することにより行う。話者方向から発せられた音響信号が入力された場合、マイクロホン51のパワー強度はマイクロホン52のそれに比べて当然大きな値となる。そこで、フレーム毎のマイクロホン51のパワー値を P_1 、マイクロホン52のパワー値を P_2 とすると式5の条件式を満たす場合に話者方向からの信号 (前方音) であると判定することができる。

【0045】

$$(P_1 - P_2) > c_1$$

【0046】ここで c_1 は予め設定された前方音検出のためのパワー差の閾値である。なお前方音の判定は式6の条件式を用いても同様の判定をすることができる。 *

$$\{1 - (P_1 / P_2)\} > c_2 \quad \text{-----式6}$$

【0048】ここで c_2 は予め設定された前方音検出のためにパワー比の閾値である。上記フレーム毎に得られた判定結果から、短時間での判定結果の切り替わりを防止するため、前方音として判定されたフレームが連続して一定フレーム数以上続いたときに前方音判定結果をオンにし、また前方音と判定されないフレームが一定フレーム数以上続いたときに前方音判定結果をオフにして、そのオン、オフの情報を外部に出力する。上記の処理により話者方向からの信号のみを検出することが可能となる。

【0049】音声判定部54では第1のマイクロホン51への入力信号が音声であるか否かを判定する。音声判定部54の動作は上記音声検出装置の第1の実施例の音声判定部3の動作と同一であるので説明は省略する。

【0050】総合判定部55では前方音検出部53、音声判定部54から一定時間間隔毎に送られてくる出力結果をもとに、各マイクロホンの組の中で話者方向からの入力が存在すると判定された第1のマイクロホンの入力信号について、音声判定部54でそれが音声信号であると判定されている場合にそのマイクロホンはオンであるという信号を外部に出力する。

【0051】以上のように本実施例によれば、話者の前後に向いた2本のマイクロホンの組を用いて、それぞれの入力信号のパワー値の違いから話者方向からの信号であるか否かを判定し、また入力信号の音韻性から音声信号であるか否かを判定するようにしたことにより、雑音による誤判定を防止し、話者方向から発せられる音声信号のみを正確に検出することができる。

【0052】次に本発明の音声検出装置の第3の実施例について図面を参照しながら説明する。図6は本実施例の動作を示すブロック図である。図6において、Wは音声を発する話者（例えば、話者W1、W2などで構成されている）、61は話者方向を向いた第1のマイクロホン（例えば、マイクロホンM11、M21などで構成されている）、62は話者と反対方向を向いた第2のマイクロホン（例えば、マイクロホンM12、M22などで構成されている）、ここで、第1のマイクロホン61と第2のマイクロホン62は、一対ごとに一組のマイクロホン（例えば、マイクロホンの組Mc1、Mc2など）として複数組のマイクロホンで構成されている。また図6において、63は第1のマイクロホンと第2のマイクロホンのそれぞれの入力信号の差から話者方向からの信号のみを検出する前方音検出部、64は各第1のマイクロ

【数5】

$$\text{-----式5}$$

* 【0047】

【数6】

$$\text{-----式6}$$

ホンの入力信号についてそのスペクトルの特徴量を検出することにより音声信号であるか否かを判定する音声判定部、65は隣合う第1のマイクロホンの組毎に入力信号間の相関をみることにより話者の位置を推定し、その話者に対応するマイクロホンを特定する話者検出部、66は上記前方音検出部63、音声判定部64、話者検出部65の出力結果をもとに最終的に各第1のマイクロホンについて前方からの音声信号が入力されているか否かを判定し、この判定結果を出力する総合判定部である。

【0053】以下、本実施例の動作を説明する。各マイクロホンに入力された音響信号はデジタル信号に変換され、全てのマイクロホン出力が前方音検出部63へ、各第1のマイクロホンの出力信号が音声判定部64、話者検出部65に送られる。

【0054】ここで前方音検出部63の動作は第2の実施例における図5の前方音検出部53の動作と同一であり、音声判定部64および話者検出部65の動作は、それぞれ第1の実施例における図1の音声判定部3、話者検出部2の動作と同一であるので説明は省略する。

【0055】総合判定部66では、前方音検出部63で前方の話者からの入力があると判定された第1のマイクロホン61が、話者検出部65でも特定された場合に、音声判定部64でその入力信号が音声であると判定されている場合に、その第1のマイクロホンはオンであるという信号を外部に出力する。

【0056】以上のように本実施例によれば、前方音検出部63で話者の前後を向いた2つのマイクロホンの組毎にその入力信号間のパワー値の違いから前方からの信号のみを検出し、音声判定部64で音韻性の検出に基づき音声信号であるか否かを判定し、話者検出部65で隣合うマイクロホンの入力信号間の相互相関係数から話者の位置を推定することにより前方からの入力のあるマイクロホンを特定し、これらの結果を総合的に判断して各マイクロホンの音声検出結果を出力するようにしたことにより、あらゆる方向からの様々な雑音が入力されても確実に棄却することができ、音声は他のマイクロホンに混入した場合でも発言した話者に対応するマイクロホンを正確に特定することができる。

【0057】次に本発明の映像切り替え装置の一実施例について図面を参照しながら説明する。図7は本実施例の構成を示すブロック図である。図7において71は各マイクロホンの入力信号からそれぞれに対応する話者の

音声信号のみを検出し、マイクロホン毎の音声信号の入力があるか否かの情報を一定時間間隔毎に出力する音声検出部、72は話者の音声が入力されているマイクロホンの位置に映像を切り換えるように制御信号を送出する映像切り替え制御部、73は、映像切り替え制御部72の出力を受けて、予め設定された発言している話者の位置にモニター74の映像を切り換えるように、カメラ75およびモニター制御部76を制御するカメラ制御部である。

【0058】以下、本実施例の動作を説明する。ここで音声検出部71は、上記で説明した音声検出装置の第1の実施例あるいは第2の実施例あるいは第3の実施例のいずれかの構成であればよく、動作の説明は省略する。

【0059】音声検出部71からは一定時間間隔毎に音声の検出されたマイクロホンの情報が出力される。この出力を受けて映像切り替え制御部72では映像切り替えのタイミングを定め、音声検出されているマイクロホン位置の映像に切り換えるよう制御信号をカメラ制御部73に送出する。ここで映像切り替えのタイミングは、映像の頻繁に切り替わることによる画面の見ずらさを回避し、また音声検出の誤検出の場合にも対応できるように、音声検出が開始されてから一定時間後に映像切り替えの信号を送出し、また音声検出が終了した時点から一定時間後に終了信号を送出する。

【0060】カメラ制御部73では、映像切り替え制御部72からの切り替え制御信号に基づき、判定されたマイクロホンに対応する話者の画面に切り換えるようにカメラ75に移動信号を送りカメラ75の向きを変更する。なお各マイクロホンに対応する話者の位置はそれぞれ予め設定しており、その位置情報がカメラ制御部73に記憶されている。

【0061】以上のように本実施例によれば、複数のマイクロホンから対応する話者の音声が入力されているもののみを正確に捉え、この音声検出情報をもとにその話者の方に自動的に映像を切り換えることが可能となり、特に自然なテレビ会議の進行を実現することのできる映像切り替え装置が実現できる。

【0062】この実施例では、一台のカメラ75を使用して、カメラ制御部73が、映像切り替え制御部72からの切り替え制御信号に基づき、判定されたマイクロホンに対応する話者に画面を切り換えるようにカメラ75に移動信号を送り、カメラ75の向きを変更するよう構成したが、複数台のカメラを、各カメラが適当数の話者に対応するように配置して、カメラ制御部73が、映像切り替え制御部72からの切り替え制御信号に基づき、判定されたマイクロホンに対応する話者に対応して配置されたカメラに接続を切り替えて、この話者に画面を切り換えるように構成することもできる。これにより、話者に対する画面の切り換えの追従性が向上して、話者の速い立ち代わりにも、十分対応できる。

【0063】

【発明の効果】請求項1の構成によれば、音声判定部が、マイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定し、話者検出部が、隣接したマイクロホンの入力信号の間の差異を検出することにより話者の位置を推定し、この話者に対応したマイクロホンを特定するので、音声判定部と話者検出部の出力結果に基づいて、総合判定部がそれぞれのマイクロホンに対応した話者の音声のみが判定できる。そのため、発声している話者に対応するマイクロホンを正確に特定することができ、様々な雑音が入力されても音声と誤検出することのない精度の高い音声検出ができる。

【0064】請求項3の構成によれば、前方音検出部が、話者方向に向いた第1のマイクロホンと話者と反対方向に向いた第2のマイクロホンに入力された信号の差異を検出して、第1のマイクロホンの前方より発せられた信号のみを検出し、音声判定部が、第1のマイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定するので、前方音検出部と音声判定部の出力結果に基づいて、総合判定部がそれぞれの第1のマイクロホンに対応した話者の音声のみが判定できる。そのため、左右、後方からの雑音、音声を棄却でき、様々な雑音が入力されても音声と誤検出することのない精度の高い音声検出ができる。

【0065】請求項4の構成によれば、前方音検出部が、一組にされた話者方向に向いた第1のマイクロホンと話者と反対方向に向いた第2のマイクロホンに入力された信号の差異を検出して、第1のマイクロホンの前方より発せられた信号のみを検出し、音声判定部が、各組の第1のマイクロホンに入力された信号からスペクトルの特徴量を抽出し、予め求めた音声の特徴量との類似性の有無によりその信号が音声であるか否かを判定し、話者検出部が、隣接した第1のマイクロホンの入力信号の間の差異を検出することにより話者の位置を推定し、この話者に対応したマイクロホンを特定するので、前方音検出部と音声判定部と話者検出部の出力結果に基づいて、総合判定部が各組の第1のマイクロホンに対応した話者の音声のみが判定できる。そのため、左右、後方からの雑音、音声を棄却でき、また発声している話者に対応するマイクロホンを正確に特定することができ、様々な雑音が入力されても音声と誤検出することのない精度の高い音声検出ができる。

【0066】請求項25の構成によれば、請求項1に記載の音声検出装置の出力に基づいて、映像切り替え制御部が、特定したマイクロホンに対応した話者に映像を切り換える制御信号をカメラ制御部に出力するので、この制御信号により、カメラ制御部が予め記憶した話者の位

置情報に基づいて出力映像の切り替えが制御できる。そのため、音声入力があったマイクロホンの位置に自動的に映像を切り換えることができ、正確で使い勝手のよい、特にテレビ会議システムでのスムーズな会議進行が実現できる。

【図面の簡単な説明】

【図1】本発明の第1の実施例の音声検出装置の構成図

【図2】同実施例の話者の特定動作の説明図

【図3】同実施例の話者の特定動作のフローチャート図

【図4】同実施例の音声判定部の構成図

【図5】本発明の第2の実施例の音声検出装置の構成図

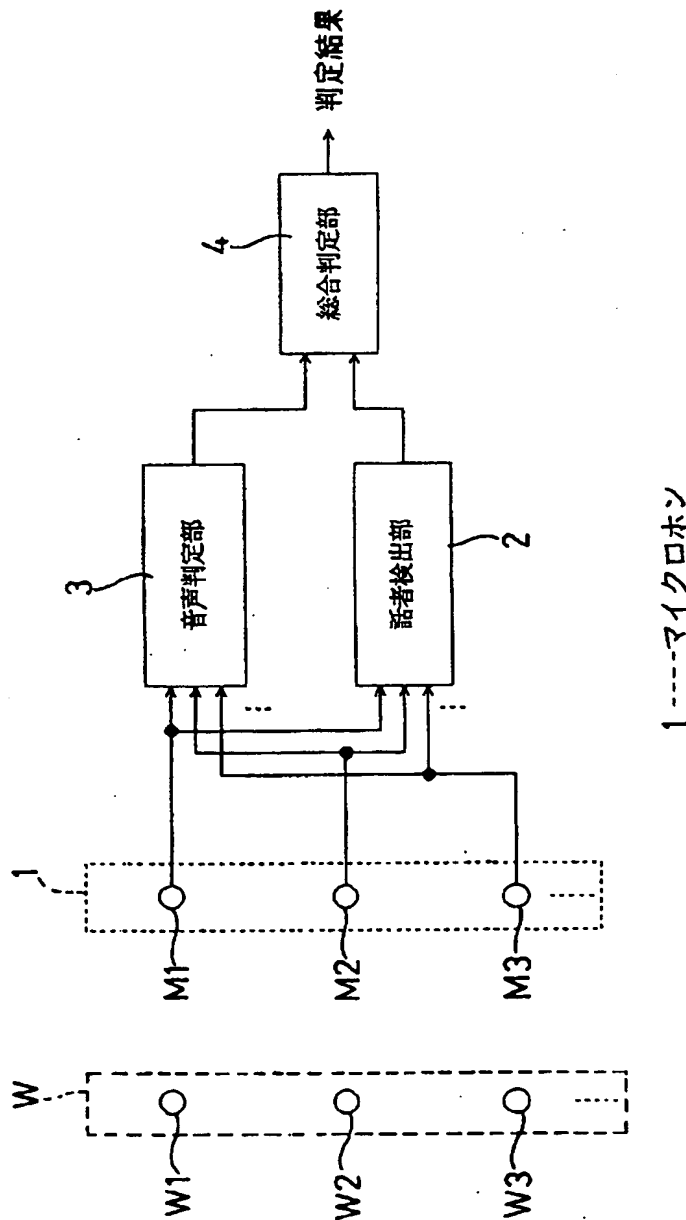
【図6】本発明の第3の実施例の音声検出装置の構成図

【図7】本発明の一実施例の映像切り替え装置の構成図

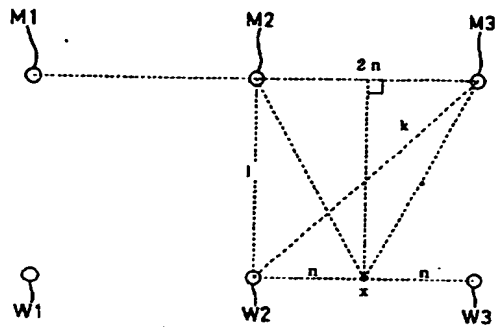
【符号の説明】

1	マイクロホン
2, 6 5	話者検出部
3, 5 4, 6 4	音声判定部
4, 5 5, 6 6	総合判定部
5 1, 6 1	第1のマイクロホン
5 2, 6 2	第2のマイクロホン
10 5 3, 6 3	前方音検出部

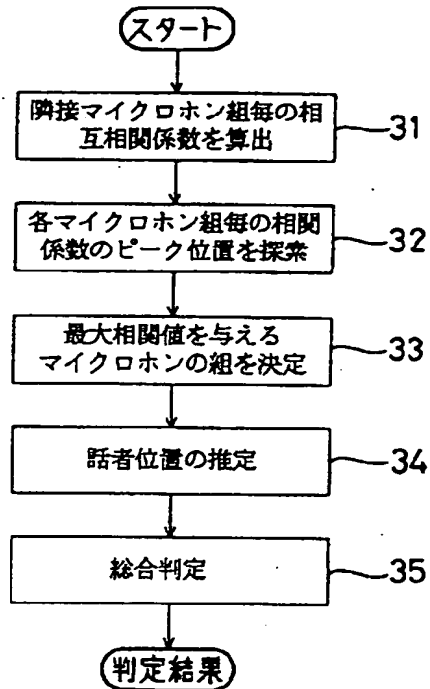
【図1】



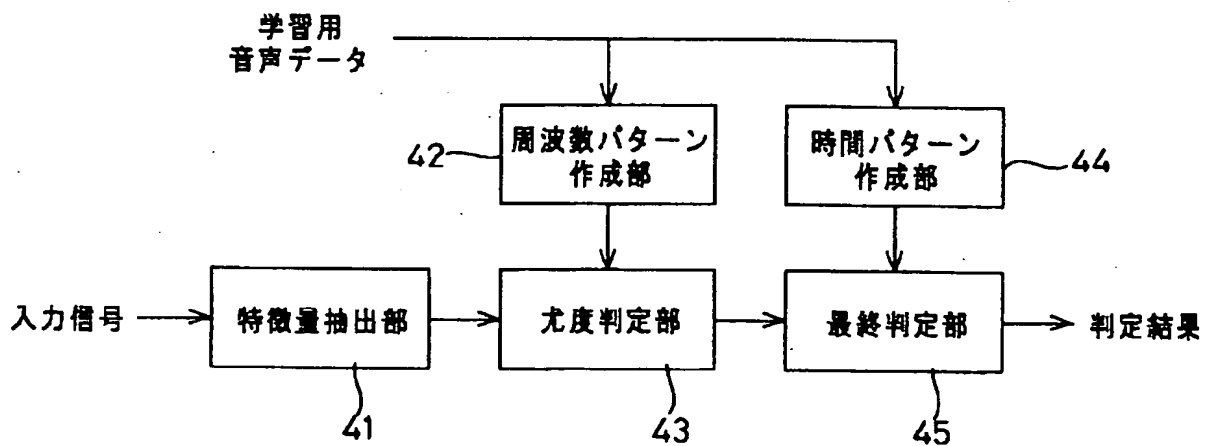
【図2】



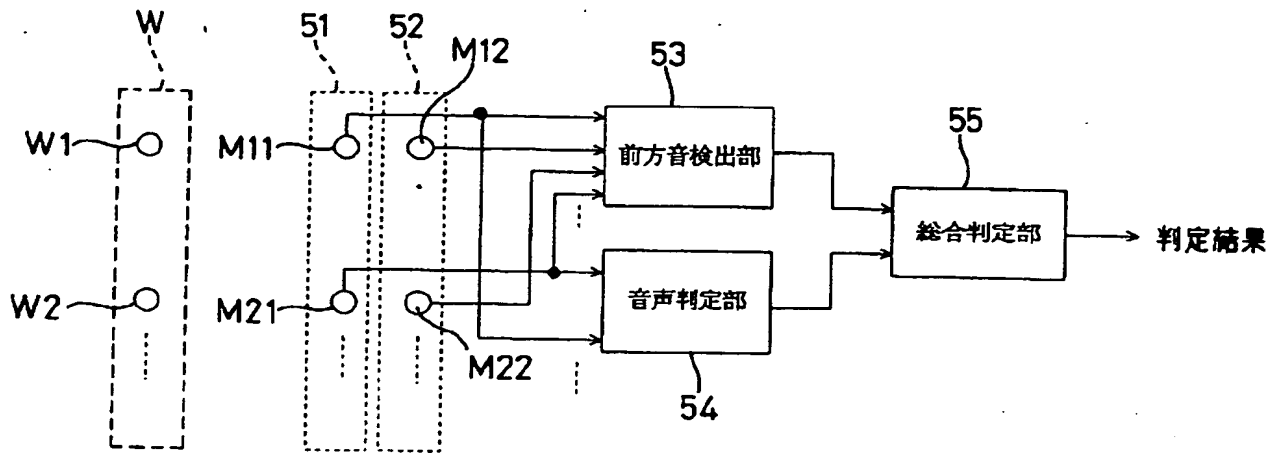
【図3】



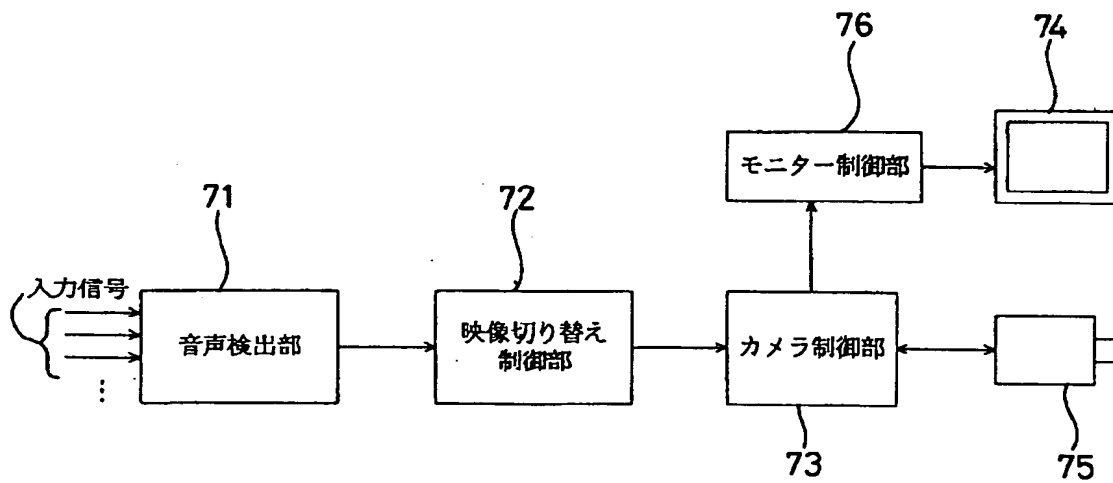
【図4】



【図 5】



【図 7】



【図6】

